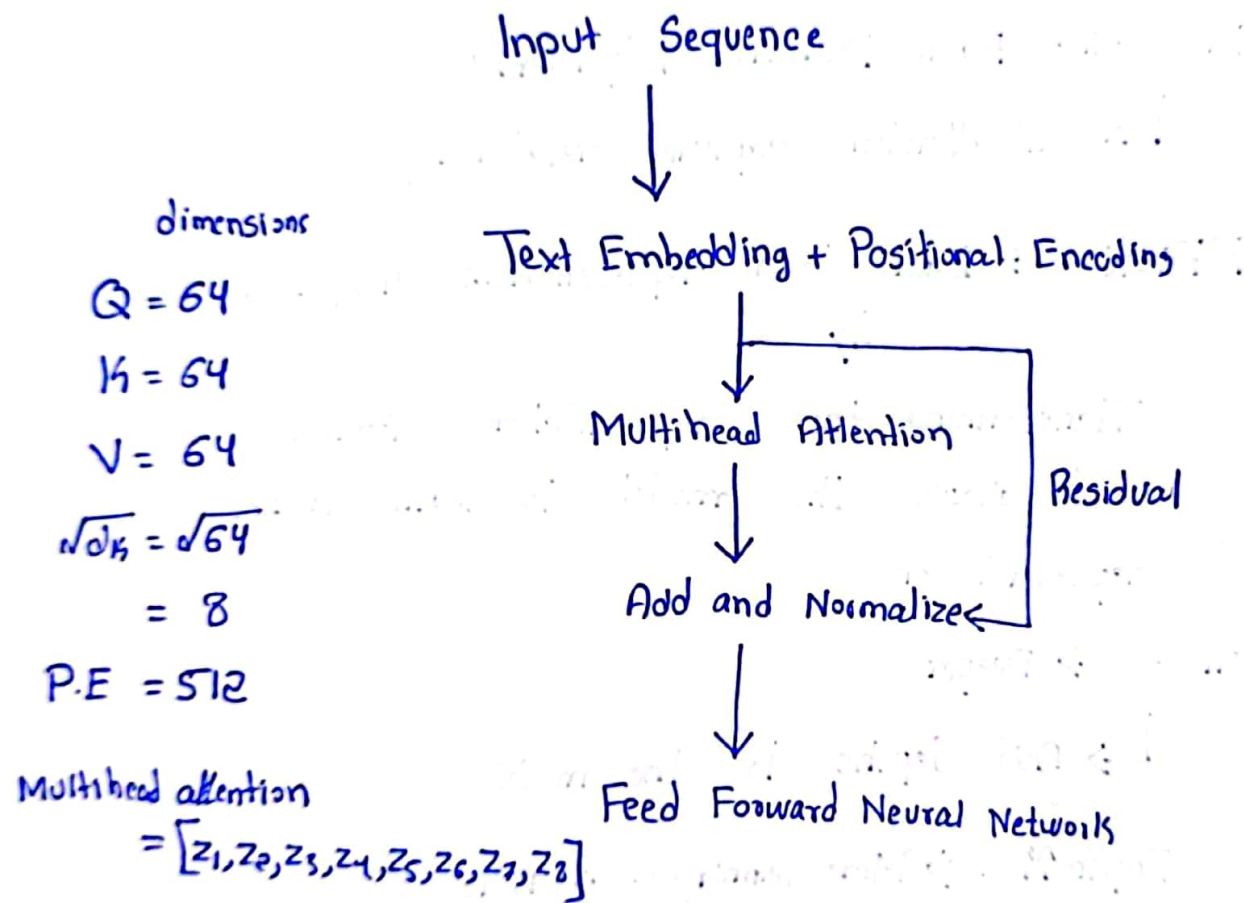


→ Complete Encoder Transformer Architecture



> Residual Connection

↳ Skip connection

① Addressing the Vanishing Gradient Problem

↳ Residual connections create a short path for gradients to flow directly through the network. Gradients remain sufficiently large.

② Improve gradient flow

↳ Convergence will be faster

③ Enables training of Deeper Networks

> Feed Forward Neural Network

① Add Non-Linearity

② Processing Each Position Independently

↳ Self attention captures relationship

FFNN → Each token representation independently

⇓

Transforming these representations further
and allows the model to learn richer
representation

③ FFNN → Deeper

↳ Add depths to the model

Depth ↑ ⇒ More learnings → DATA